

DUM: Diversity-Weighted Utility Maximization for Recommendations

Azin Ashkan *, Branislav Kveton *, Shlomo Berkovsky **, Zheng Wen ***

**Technicolor, United States*

{azin.ashkan, branislav.kveton}@technicolor.com

***CSIRO, Australia*

shlomo.berkovsky@csiro.au

****Yahoo Labs, United States*

zhengwen@yahoo-inc.com

ABSTRACT

The need for diversification of recommendation lists manifests in a number of recommender systems use cases. However, an increase in diversity may undermine the utility of the recommendations, as relevant items in the list may be replaced by more diverse ones. In this work we propose a novel method for maximizing the utility of the recommended items subject to the diversity of user's tastes, and show that an optimal solution to this problem can be found greedily. We evaluate the proposed method in two online user studies as well as in an offline analysis incorporating a number of evaluation metrics. The results of evaluations show the superiority of our method over a number of baselines.

Keywords

Recommender systems, polymatroid, diversity, utility

1. INTRODUCTION

The popularity of recommender systems has soared in the recent years. They are widely used in social networks, entertainment, eCommerce, Web search, and many other online services [20]. Recommenders deal with the information overload problem and select items on behalf of their users. Typically, a recommender scores recommendable items according to their match to the user's preferences and interests, as encapsulated in the user profiles, and then recommends a list of top-scoring items.

A naïve selection of top-scoring items may, however, yield a sub-optimal recommendation list. For instance, collaborative recommenders tend to recommend most popular items appearing in the profiles of numerous users [13]. While being good recommendations on their own, these items are likely to be known to the user and bear little value. Likewise, content-based recommenders may target user's favorite topics and recommend homogeneous lists that overlook other potentially interesting topics [16]. This has brought to the fore the problem of diversity in recommender systems, which has been studied in a number of works [5, 10, 17, 24, 31].

In a nutshell, the diversity problem deals with the construction of recommendation lists that cover as wide as possible range of topics of interest. The problem is particularly acute for users with eclectic interests, having no single dominant topic but rather interested in multiple topics. In this case, it is important for the recommendation list to include items that touch upon many topics, in order to increase the chance of answering the current user's need. Repercussions of the diversity problem can be recognized also in other recommender system use cases. Consider the group recommendation problem in heterogeneous (in terms of interests) groups. Another example of the need for diversity is in sequential recommendations, like in the music domain. In both cases, the recommendation list should incorporate diverse items that either appeal to a number of group members or represent a number of music genres [29].

The need for diversity manifests itself also beyond recommender systems. Consider an ambiguous Web search query. Having no knowledge about the context of the query, a search engine may present results pertaining to different interpretations of the query, so that the user can pick the desired one and reformulate the query [3]. Another instance comes from text summarization. Unless the desired topic of the summary is known, it should include references to as many aspects of the original document as possible. Also, diversification may be useful in computer supported collaborative work. There, formation of virtual groups may need to bring together users with complementary skills and expertise areas, such that the diversity of the group is important.

In all the above diversification use cases, it is of paramount importance to maintain the trade-off between increasing the list diversity and maintaining the utility of the results [4, 11, 30]. Diversity typically comes at the account of decreasing the relevance of items, as relevant but redundant items are substituted with less relevant but more diverse ones. Hence, there is a need to strike the balance between the two objectives [4], a modular relevance function and a submodular diversity function, for which an approximation to the optimal solution can be computed greedily [18].

In this work, we introduce a different objective diversification function and show that an optimal solution to the diversity problem can be found greedily. We propose a parameter-free method, denoted as *diversity-weighted utility maximization* (DUM), which maximizes the utility of the items recommended to users, subject to the diversity of their interests. We cast this problem as finding the maximum of a modular function subject to a submodular constraint [8], which is known to have an optimal greedy solution. This solution guarantees that items in the recommendation list cover different interests in the user profile, such that each topic of interest is represented by items with high utility. In other words, the utility of items

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

remains the primary concern, but it is subjective to maintaining the diversity and avoiding redundancy in the list. We discuss several interpretations of DUM and identify suitable submodular diversity functions.

We conduct an extensive evaluation of the proposed approach. We present two online studies using crowdsourcing, which compare the perceived quality of the lists generated by DUM with baseline settings maximizing a linear combination of utility and diversity. The results show the superiority of the lists generated by DUM over the baseline methods and we characterize the cases when this superiority is prominent. We also present an offline evaluation that applies a variety of metrics to (a) exemplify the trade-off between diversity and utility in recommendations; and (b) demonstrate that DUM successfully outperforms the baselines. Overall, our analyses demonstrate that DUM can effectively deliver personalized recommendations with high degree of utility and diversity, while not requiring a-priori parameter tuning.

In summary, the contribution of this work is two-fold. Firstly, we propose a parameter-free and computationally efficient method aimed at improving the diversity of the recommendation lists, while maintaining their utility. Secondly, we present experimental evaluations – online user studies and offline experiments alike – that demonstrate solid empirical evidence supporting the validity of the proposed approach.

Note: The following notation is used throughout this paper. Let A and B be sets, and e be an element of a set. We use $A + e$ instead of $A \cup \{e\}$ and $A + B$ instead of $A \cup B$. Furthermore, we use $A - e$ instead of $A \setminus \{e\}$ and $A - B$ instead of $A \setminus B$. We represent *ordered sets* by vectors and also refer to them as *lists*.

2. RELATED WORK

A common approximation to diversified ranking is based on the notion of *maximal marginal relevance* (MMR) proposed by Carbonell and Goldstein [4]. In this approach, utility (e.g., relevance) and diversity are represented by independent metrics. Marginal relevance is defined as a weighted combination of these two metrics, to account for the trade-off between utility and diversity. Given a standard ranking of items, R , a diversified re-ranking of these items, S , is created (Algorithm 1). In each iteration, an item $e^* \in R - S$ is chosen, such that it maximizes the marginal relevance:

$$e^* = \arg \max_{e \in R - S} (1 - \lambda) \mathbf{w}(e) + \lambda f(S + e) \quad (1)$$

where $\mathbf{w}(\cdot)$ and $f(\cdot)$ represent the notions of utility and diversity, respectively, and the parameter λ controls the trade-off between the two. Typically, the utility \mathbf{w} is a modular function of S , whereas the diversity f is a submodular function of S . The existing approaches differ in how they account for different aspects of query or user (or any other entity of interest) to model $f(\cdot)$.

Implicit approaches assume that similar items should be penalized since they cover similar aspects. For instance, Yu et al. [26] compute $f(S + e) = -\max_{e' \in S} \text{Sim}(e, e')$ to measure the redundancy of user intent e with respect to a set of selected intents S , where $\text{Sim}(e, e')$ is the cosine similarity between the user intents e and e' . Gollapudi and Sharma [9] propose multiple diversification objectives considering the tradeoff between relevance and diversity and using various axioms, relevance functions, and distance functions. Their distance functions are defined based on various implicit metrics, e.g., document content, in order to capture the pairwise similarity between any pair of documents.

On the other hand, explicit approaches model different aspects (e.g., query intent, query topic, or movie genre) directly, and promote diversity by maximizing the coverage of selected items with

Algorithm 1 MMR: Maximal Marginal Relevance

Input:

Standard ordering of items R

$S \leftarrow (), n = |R|$

while $|S| < n$ **do**

$e^* \leftarrow \arg \max_{e \in R - S} \lambda \mathbf{w}(e) + (1 - \lambda) f(S + e)$

$R \leftarrow R - e^*$

Append item e^* to list S

end while

Output:

List of recommended items S

respect to these aspects. For instance, Santos et al. [21] define $f(S + e) = \sum_{t \in \mathcal{T}_q} P(t|q)P(e, \tilde{S}|t)$ where $P(d, \tilde{S}|t)$ represents the likelihood of document e satisfying topic t while the ones in S fail to do so. Also, $P(t|q)$ denotes the popularity of t among all possible topics \mathcal{T}_q that may satisfy a user’s information need from issuing query q .

In addition to the above approaches in diversifying existing rankings, another group of work directly learns a diverse ranking by maximizing a submodular objective function. Among these approaches, Radlinski et al. [19] and Yue and Guestrin [27] propose *online* learning algorithms for optimizing a class of submodular objective functions for diversified retrieval and recommendation, respectively. Agrawal et al. [1], on the other hand, address search result diversification in an offline setting, with respect to the topical categories of documents. The authors target the maximization of a submodular objective function following the definition of marginal relevance. They propose a greedy algorithm to approximate the objective function and show that an optimal solution can be found in a special case, where each document belongs to exactly one category.

Vallet and Castells [23] study personalization in combination with diversity such that the two objectives complement each other in addressing various query aspects and satisfying user needs. In particular, they generalize the work of Agrawal et al. [1] and Santos et al. [21] to the personalized versions by exploiting available information about user preferences.

Most of these studies target diversity in information retrieval, while there has been a growing interest in recommendation diversification more recently. One of the initial works in recommendation diversification is by Ziegler et al. [31], who argue that user satisfaction does not solely depend on the accuracy of recommendation results. The authors propose a similarity metric, the intra-list similarity (ILS), which computes the average pairwise similarity of items in a list. A higher value of the metric denotes a lower diversity. They use this metric in their topic diversification model to control a balance between the accuracy and diversity of recommendations.

Zhang et al. [28] formulate the diversification problem as finding the best possible subset of items to be recommended over all possible subsets. They address this as the maximization of the diversity of a list of recommended items, subject to maintaining the accuracy of the items. Zhou et al. [30] propose a hybrid method that targets the maximization of a weighted combination of independent utility- and diversity-based approaches requiring parameter tuning to control the tradeoff between the objectives of the two approaches.

Most of the existing diversification approaches consider the maximization of an objective function to satisfy a user’s need in terms of utility and diversity of the result list. Most of these approaches are based on the idea behind the maximal marginal relevance where a submodular objective function (Equation 1) is maximized. There-

fore, a $(1 - 1/e)$ -approximation to the optimal solution can be computed greedily [18]. This paper is an extension to our prior work in [2], where we introduce a new objective function for recommendation diversification, the optimal solution of which can be found greedily. This objective function targets the utility as the primary concern, and maximizes it subject to maintaining the diversity of user's tastes. In this paper, we elaborate on the intuitions and details behind the proposed greedy algorithm, and provide extensive online and offline evaluations on its performance in practice. We show that this method is computationally efficient and parameter-free, and it guarantees that high-utility items appear at the top of the recommendation list, as long as they contribute to the diversity of the list.

3. MOTIVATING EXAMPLES

In this section, we discuss several motivating examples for our work. A more formal description of our method and its analysis are presented in Section 4.

Consider the following recommendation problem. A system recommends to a user movies from a ground set of four movies:

ID e	Movie name	Utility $w(e)$	Action	Comedy
1	Inception	0.8	X	
2	Spider-Man 2	0.7	X	
3	Grown Ups 2	0.5		X
4	The Sweep	0.2		X

The user likes either *Action* or *Comedy* movies, depending on the mood of the user, but the system does not know the user's mood. The user chooses the first recommended movie e that matches the genre that the user currently prefers and is satisfied proportionally to the utility of the movie $w(e)$, the probability that e is liked. Our goal is to recommend a minimal list of movies that maximizes the user's satisfaction and also covers all user's preferences, irrespective of the user's mood.

The optimal solution to our problem is a list of two movies, $S = (1, 3)$. When the user prefers *Action* movies, the user selects the first recommended movie in the list, *Inception*, and is satisfied with probability 0.8. This is substantially greater than if *Spider-Man 2*, another *Action* movie, was in the list instead of *Inception*. On the other hand, when the user prefers *Comedy* movies, the user selects the second recommended movie in the list, *Grown Ups 2*, and is satisfied with probability 0.5. This is substantially greater than if *The Sweep*, another *Comedy* movie, was in the list instead of *Grown Ups 2*. Note that the solution S can be computed greedily. In particular, S is a list of two highest-utility movies, one from each genre.

Now suppose that we add to the ground set a movie that is both *Action* and *Comedy*, and its utility is 0.7:

ID e	Movie name	Utility $w(e)$	Action	Comedy
1	Inception	0.8	X	
2	Spider-Man 2	0.7	X	
3	Grown Ups 2	0.5		X
4	The Sweep	0.2		X
5	Kindergarten Cop	0.6	X	X

The optimal solution to the problem is a list $S = (1, 5)$. When the user prefers *Action* movies, the user selects the first recommended movie, *Inception*, and is satisfied with probability 0.8. On the other hand, when the user prefers *Comedy* movies, the user selects the second recommended movie, *Kindergarten Cop*, and is satisfied

with probability 0.6. Note again that the solution S can be computed greedily. It is a list of two highest-utility movies, one from each genre.

Finally, we replace the last movie in the ground set with a movie whose utility is 0.9:

ID e	Movie name	Utility $w(e)$	Action	Comedy
1	Inception	0.8	X	
2	Spider-Man 2	0.7	X	
3	Grown Ups 2	0.5		X
4	The Sweep	0.2		X
5	Indiana Jones and the Last Crusade	0.9	X	X

The optimal solution to the problem is a single movie, $S = (5)$. The reason is that *Indiana Jones and the Last Crusade* is the highest-utility movie in both *Action* and *Comedy*. Hence, it is the best recommendation irrespective of the user's mood. Note again that the solution S can be computed greedily. It is the highest-utility movie that belongs to both genres

In all three examples, the optimal solutions can be computed greedily. This is not by chance. In the next section, we generalize the ideas expressed in these examples and introduce the notion of diverse recommendations where the optimal solution can be found greedily. This is the main contribution of our paper.

4. DIVERSITY-WEIGHTED UTILITY MAXIMIZATION

Our objective is to maximize the utility of recommending a list of items to a user subject to the diversity of their tastes. We present the formal definition of our method in Section 4.1, followed by Section 4.2 where we show that the optimal solution of the method can be found efficiently. The interpretations and intuitions behind our method are explained in Section 4.3. We show that the length of the list recommended by our method can be controlled by considering different diversity constraints dependent on user preferences in Section 4.4.

4.1 Problem Formulation

Let $E = \{1, \dots, L\}$ be a ground set of L recommendable items, such as movies or songs. Let $\mathbf{w} \in (\mathbb{R}^+)^L$ be a vector of item utilities, such as item popularity scores or predicted ratings. The e -th entry of \mathbf{w} , $w(e)$, is the utility of item e .

The objective of the diversification method is to maximize the satisfaction of the user subject to the diversity of their tastes. However, an increase in diversity typically comes at the account of a decrease in the utility of the items in the list, e.g., relevant but redundant items are substituted by less relevant but more diverse items. Addressing this tradeoff and striking the balance between increasing the diversity and maintaining the utility is an important challenge for any diversification method. Considering the *utility* as the primary concern, we aim to expose the user to a variety of *choices* in the recommendation list, while losing the minimal amount of utility in the provision of these choices.

In order to recommend a list of items that maximizes the user's utility of choice, we target at maximizing the utility of the recommendation list weighted by the increase in diversity. In other words, each increase in diversity is covered by the item with the highest possible utility. Formally, our diversity-weighted utility maximiza-

tion (DUM) problem is formulated as:

$$A^* = \arg \max_{A \in \Theta} \sum_{k=1}^L g_A(a_k) \mathbf{w}(a_k), \quad (2)$$

where $A = (a_1, \dots, a_L)$ is an ordered set of items E , Θ is the set of all permutations of E , and $A^* = (a_1^*, \dots, a_L^*)$ is the optimal solution to the problem. The vector $g_A \in (\mathbb{R}^+)^L$ are the gains in diversity associated with items E . In particular:

$$g_A(e) = f(A_{k-1} + e) - f(A_{k-1}) \quad (3)$$

is the gain in diversity associated with choosing item e given a set of previously chosen items in A , where k is such that $a_k = e$ and $A_k = \{a_1, \dots, a_k\}$ is an unordered set of the first k items in A . The function $f : 2^E \rightarrow \mathbb{R}^+$ is a diversity function from subsets of the ground set E to non-negative real numbers.

The diversity function f can have many different forms. For instance, $f(X)$ can be the number of unique genres covered by movies X recommended by a recommender system. Alternatively, $f(X)$ can be the average pairwise dissimilarity between a set of products X recommended by a shopping website. In this work, we assume that the diversity function f is *monotonically increasing*:

$$\forall X \subseteq E, e \in E - X : f(X + e) - f(X) \geq 0, \quad (4)$$

the diversity of any set X does not decrease when any item e is added to this set. This assumption is quite natural. We also assume that $f(\emptyset) = 0$, the diversity of the empty set is zero. This assumption is without loss of generality. In particular, it can be always satisfied by subtracting $f(\emptyset)$ from f .

4.2 Greedy Solution

For a general monotonic function f , the optimization problem (2) is NP-hard. However, when f is submodular, the problem can be cast as finding a maximum-weight basis of a polymatroid [8] and can be solved greedily. We first present the greedy algorithm and then argue that it is optimal.

The pseudo-code of the greedy algorithm for *diversity-weighted utility maximization* (DUM) is shown in Algorithm 2. The algorithm works as follows. First, the items E are sorted in decreasing order according to their utility, $\mathbf{w}(a_1^*) \geq \dots \geq \mathbf{w}(a_L^*)$, and placed into $A^* = (a_1^*, \dots, a_L^*)$. Then we examine the items in this order. When $g_{A^*}(a_k^*) > 0$, item a_k^* is added to the list of recommended items S . When $g_{A^*}(a_k^*) = 0$, item a_k^* is not added to S because it does not contribute to the diversity of S . Finally, the algorithm returns the recommendation list S .

We illustrate DUM on the second example in Section 3. In this example, $A^* = (1, 2, 5, 3, 4)$, and the diversity gains of movies 2, 3 and 4 are zero due to the contribution of their preceding movies in the list. Therefore, these movies are not placed into the recommendation list, and $S = (1, 5)$.

DUM has several notable properties. First, it is *parameter-free*. That is, DUM does not require any parameter tuning and therefore should be robust in practice. Second, DUM is a greedy method and therefore is *computationally efficient*. In particular, suppose that the diversity function f is an oracle that can be queried in $O(1)$ time. Then the time complexity of DUM is $O(L \log L)$, comparable to the complexity of sorting L numbers. Finally, DUM computes the *optimal solution* to the optimization problem (2).

In the rest of Section 4, we analyze DUM both in terms of A^* and S . Note that the solutions A^* and S are equivalent in the sense that S is a list obtained from A^* by eliminating the items that have zero contribution in the objective function (2). Therefore, the values of

Algorithm 2 DUM: Diversity-Weighted Utility Maximization

Input:

Ground set E
Weight vector \mathbf{w}

// Compute the maximum-weight basis of a polymatroid

Let a_1^*, \dots, a_L^* be an ordering of items E such that:

$\mathbf{w}(a_1^*) \geq \dots \geq \mathbf{w}(a_L^*)$
 $A^* \leftarrow (a_1^*, \dots, a_L^*)$

// Generate the list of recommended items S

$S \leftarrow ()$

for $k = 1, \dots, L$ **do**

$g_{A^*}(a_k^*) \leftarrow f(A_{k-1}^* + a_k^*) - f(A_{k-1}^*)$

if $(g_{A^*}(a_k^*) > 0)$ **then**

Append item a_k^* to list S

end if

end for

Output:

List of recommended items S

the solutions are identical. So the difference in treatment is purely technical and allows us to reduce overhead in notation.

The optimality of DUM can be argued based on the following observation. Our optimization problem (2) is equivalent to maximizing a modular function on a polymatroid [8], a well-known combinatorial optimization problem that can be solved greedily. In particular, let $M = (E, f)$ be a polymatroid, where E is its ground set and f is a submodular diversity function. Let:

$$P_M = \left\{ \mathbf{x} : \mathbf{x} \in \mathbb{R}^L, \mathbf{x} \geq 0, \forall X \subseteq E : \sum_{e \in X} \mathbf{x}(e) \leq f(X) \right\} \quad (5)$$

be the independence polyhedron associated with function f . Then the maximum-weight basis of M is defined as:

$$\mathbf{x}^* = \arg \max_{\mathbf{x} \in P_M} \langle \mathbf{w}, \mathbf{x} \rangle, \quad (6)$$

where $\mathbf{w} \in (\mathbb{R}^+)^L$ is a vector of non-negative weights. Because P_M is a submodular polytope and the weights \mathbf{w} are non-negative, the optimization problem (6) is equivalent to finding the order of dimensions A in which $\langle \mathbf{w}, \mathbf{x} \rangle$ is maximized [8]. This problem can be written formally as (2) and has the same greedy solution as in DUM. In particular, the items E are sorted in decreasing order according to their weights, $\mathbf{w}(a_1^*) \geq \dots \geq \mathbf{w}(a_L^*)$, and placed into $A^* = (a_1^*, \dots, a_L^*)$. Finally, $\mathbf{x}^* = g_{A^*}$.

4.3 Interpretation

In this section, we discuss several interpretations of DUM. Without loss of generality, we assume that the different aspects of user's taste are represented by a finite set of topics $\mathcal{T} = \{1, \dots, M\}$. For example, in a movie recommendation domain, these topics can be the genres of movies, such as $\mathcal{T} = \{\text{Drama}, \text{Comedy}, \text{Action}\}$.

Our first observation is that if the diversity of a set of items is measured by the number of unique topics covered by the items, then DUM generates a list of items, where each topic is covered by the highest-utility item in that topic.

LEMMA 1. *Let the diversity function f be defined as the num-*

ber of topics covered by items X :

$$f(X) = \sum_{t \in \mathcal{T}} \mathbb{1}\{\exists e \in X : \text{item } e \text{ covers topic } t\}.$$

Then DUM returns a recommendation list S , where each topic t is covered by the highest-utility item that belongs to t . Moreover, the length of S is at most $|\mathcal{T}|$.

PROOF. The first claim is proved by contradiction. Let e_t^* be the item with the highest utility that belongs to topic t . Suppose that item e_t^* is not chosen by DUM, e_t^* is not in list S generated by DUM. Then $g_{A^*}(e_t^*) = 0$, which implies that another item must have covered topic t before item e_t^* . However, this is a contradiction, since e_t^* is the item with the highest utility from t , and therefore DUM must have tested it before any other item that covers t .

The second claim follows from the fact that $g_{A^*}(a_k^*) > 0$ implies that the value of the diversity function f increases by at least one. By definition, $f(X) \leq |\mathcal{T}|$ for any X . Therefore, the maximum number of items added to S is $|\mathcal{T}|$. ■

Our second observation is that our objective (2) can be viewed as maximizing the *expected utility* of choosing an item when the diversity gains $g_A(e)$ are viewed as the probabilities of choosing items. This interpretation is motivated by the cascade model [7] of user behavior, which considers the relationship between successive items in a list. In this model, users scan the list from top to the bottom and eventually stop because either their information need is satisfied or their patience is exhausted.

Specifically, note that for any ordering of items A :

$$\begin{aligned} \sum_{k=1}^L g_A(a_k) &= \sum_{k=1}^L [f(A_{k-1} + a_k) - f(A_{k-1})] \\ &= f(E) - f(\emptyset) + \sum_{k=1}^{L-1} [f(A_k) - f(A_{k-1})] \\ &= f(E). \end{aligned} \quad (7)$$

The first equality is due to the definition of the diversity gains (3). The second equality follows from the fact that $A_k = A_{k-1} + a_k$. The last equality is due to the observation that $f(\emptyset) = 0$. It follows that:

$$\forall e \in E : \frac{g_A(e)}{f(E)} \in [0, 1], \quad \frac{1}{f(E)} \sum_{k=1}^L g_A(a_k) = 1, \quad (8)$$

and therefore $g_A(e)/f(E)$ can be interpreted as the probability of choosing item e , given that none of the earlier recommended items A_{k-1} is chosen. Under this assumption, $\sum_{k=1}^L g_A(a_k) \mathbf{w}(a_k)$ is the expected utility of choosing an item, scaled up by $f(E)$.

4.4 Diversity Function

The length of the recommendation list S generated by DUM depends on the diversity function f . In extreme cases, this list may include all items. For instance, consider a problem where:

$$\forall X \subseteq E, e \in E - X : f(X + e) - f(X) > 0, \quad (9)$$

the diversity increases when any item e is added to any subset of items X . Then for any ordering A , $g_A(e) > 0$ for all items e . As a result, $g_{A^*}(e) > 0$ for all items e and DUM returns the list of all items sorted in the descending order of utility. This result is mathematically correct. But it is not a very useful diverse recommendation.

To get useful diverse recommendations, it is important to control the maximum number of items returned by DUM, e.g., by choosing

appropriate diversity functions. For instance, for the diversity function in Lemma 1, the maximum length of the recommendation list S is equal to the number of topics $|\mathcal{T}|$. In this section, we generalize the ideas from Section 4.3 and propose another class of diversity functions that are suitable for DUM.

Consider the case where different users may have different tolerance for redundancy in the recommendation list due to their interests and preferences [14, 25]. These differences can be modeled by a diversity function that assigns different weights to each topic of interest. In particular, the function can be defined as:

$$f(X) = \sum_{t \in \mathcal{T}} \min \left\{ \sum_{e \in X} \mathbb{1}\{\text{item } e \text{ covers topic } t\}, N_t \right\}, \quad (10)$$

where N_t is the number of items from topic t that is required to be in the recommendation list of a given user. In the next lemma, we characterize the output of DUM for the above function.

LEMMA 2. Let the diversity function f be defined as in (10). Then DUM returns a recommendation list S such that each topic t is covered by at least N_t items of the highest utility that cover topic t . Moreover, the length of S is at most $\sum_{t \in \mathcal{T}} N_t$.

PROOF. The first claim is proved by contradiction. Let $e_{t,k}^*$ be the k -th item with the highest utility from topic t , where $k \leq N_t$. Suppose that item $e_{t,k}^*$ is not chosen by DUM, $e_{t,k}^*$ is not in list S generated by DUM. Then $g_{A^*}(e_{t,k}^*) = 0$, which implies that topic t must have been covered at least N_t times before DUM tests item $e_{t,k}^*$. However, note that this is a contradiction, since $e_{t,k}^*$ is among the first N_t items that cover topic t , and therefore among the first N_t items from that topic that are tested by DUM.

The second claim follows from the fact that $g_{A^*}(a_k^*) > 0$ implies that the value of the diversity function f increases by at least one. By definition, $f(X) \leq \sum_{t \in \mathcal{T}} N_t$ for any X . Therefore, the maximum number of items added to S is $\sum_{t \in \mathcal{T}} N_t$. ■

The diversity function in (10) allows for controlling the length of the recommendation list S . In particular, if topic t is irrelevant for the user, N_t should be set to 0. As a rule of thumb, more relevant topics t should be assigned higher weights N_t .

5. EXPERIMENTS

The proposed method is evaluated in two online user studies and in an offline evaluation. In each case, we compare the performance of DUM to variants of MMR, since many existing diversification approaches are based on the objective function of MMR (Section 2). In theory, the optimal solution to DUM can be found greedily, while MMR finds only a $(1 - 1/e)$ -approximation to the optimal solution. Through the empirical evaluation, we show that DUM satisfies the users' needs better than MMR, and it is superior in recommending lists that satisfy utility and diversity at the same time.

We conduct two online studies using Amazon's Mechanical Turk¹ (MT). In the first study, we *evaluate* separately the recommendation lists generated by DUM and MMR, by asking MT workers to identify in the lists a movie that matches their genre of interest and indicate the relevance of this movie. In the second study, we *compare* the DUM and MMR recommendation lists, by asking MT workers to judge the coverage of two movie genres by the lists. We also report the findings of an offline study, where we perform a fine-grained assessment of the evaluated methods by creating user preference profiles and considering various combinations of genres.

¹<http://www.mturk.com>

Instructions

Imagine the following situation. You want to watch a movie from a particular movie genre. The recommender system does not know this genre and recommends you a list of movies. Please choose the movie genre that you would like to watch and evaluate four lists of recommended movies. Choose the most appropriate judgments.

Choose the movie genre that you would like to watch:

drama

List 1

1. So Close
2. The Shawshank Redemption
3. Fight Club
4. The Lord of the Rings: The Fellowship of the Ring
5. Inception
6. The Lord of the Rings: The Return of the King

Which of the above movies matches your chosen genre? Answer "None" if no movie does.

None

Is this movie a good recommendation for your chosen genre? Do not answer if you answered "None" above.

Not applicable

List 2

1. The Shawshank Redemption
2. The Lord of the Rings: The Fellowship of the Ring
3. Inception
4. Titanic
5. Finding Nemo
6. The Shining

Which of the above movies matches your chosen genre? Answer "None" if no movie does.

None

Is this movie a good recommendation for your chosen genre? Do not answer if you answered "None" above.

Not applicable

Figure 1: A portion of our Mechanical Turk questionnaire in user study 1. We only show the first two lists of recommended movies.

5.1 User Study 1

In the first study, we evaluate the diversity and utility of DUM in a movie recommendation application. We compare DUM to three variants of MMR, which are parametrized by $\lambda \in \{\frac{1}{3}, \frac{2}{3}, 0.99\}$.

The ground set E are 10k most frequently rated IMDb² movies. The utility of movie e , $w(e)$, is the number of ratings assigned to this movie. The values of $w(e)$ are normalized such that the maximum utility is 1, i.e., $\max_{e \in E} w(e) = 1$. The diversity function f is defined as in Lemma 1. We also normalize f such that the maximum diversity is 1, i.e., $\max_{e \in E} f(e) = 1$. The set of topics \mathcal{T} includes 8 most popular movie genres in E :

$$\mathcal{T} = \{\text{Drama, Comedy, Thriller, Romance, Action, Crime, Adventure, Horror}\}. \quad (11)$$

We restrict our attention to 8 genres only since 8 genres can be always covered by 8 movies, a reasonably short list of movies that can be evaluated by a MT worker.

All methods are evaluated in 200 MT human intelligence tasks (HITs). In each HIT, we initially ask the worker to choose a genre of interest. Then, we generate four recommendation lists: one by DUM and three by MMR for different values of λ . The generation of the lists is independent of the genre chosen by the worker. Finally, we ask the worker to evaluate the lists. For each list, we ask two questions. First, we ask the worker to identify a movie in the list that matches the chosen genre. This question addresses the diversity of the list, whether the chosen genre is covered by the list. The worker can also answer “none” if the list does not contain a

	DUM	$\lambda = \frac{1}{3}$	$\lambda = \frac{2}{3}$	$\lambda = 0.99$
List includes a movie that matches the chosen genre	84.0%	70.5%	67.0%	66.5%
The chosen movie is a good recommendation	77.0%	64.5%	62.5%	62.5%

Table 1: Comparison of DUM and MMR in user study 1. For each method, we report the percentage of times that the worker finds a matching movie in the list and the percentage of times that the matching movie is a good recommendation.

movie from the chosen genre. If the worker identifies a matching movie, we ask the worker if the movie is a good recommendation for the chosen genre. This question addresses the utility of the list, whether the chosen genre is covered by a good movie in the list. A screenshot of our MT questionnaire is shown in Figure 1.

In each HIT, we present the four recommendation lists in a random order. This eliminates the *position bias*. In addition, in each HIT the set of recommendable movies contains 3.3k movies chosen at random from the 10k movies in E . Hence, the recommendation lists differ across the HITs, which eliminates the *item bias*, i.e., the workers cannot prefer one method over another only because the recommended movies are inherently more likable. All the recommendation lists are of the same length – the length of the list produced by DUM. We adopt this methodology because we want to compare the lists for the same number of movies in the lists. Note that that we do not put MMR into any disadvantage. In particular, for any DUM list, MMR can generate lists that are either of higher utility or more diverse than the DUM list, when the value of λ is large or small, respectively. This can be seen in Table 2, for instance.

Our HITs are completed by 34 *master* workers, who are MT’s elite workers chosen based on the high quality of the work. Each worker is asked to complete at most 8 HITs. This guarantees that our HITs are completed by more than just a handful of workers. On average, a worker spends 72 seconds per HIT, i.e., 19 seconds to evaluate a list of up to 8 movies. Later in this section, we present two permutation tests that show that our results are highly unlikely under the hypothesis that the workers are of low quality, or that the questions are answered randomly. This implies that the workers have reasonable expertise in evaluating the HITs.

The results of our study are presented in Table 1. For each compared method, we report the percentage of times that the worker finds a movie in the list that matches the chosen genre and the percentage of times that the matching movie is a good recommendation. We observe two major trends.

Firstly, the percentage of times that the worker finds a matching movie in the DUM list is 13.5% higher than in the list generated by the best performing baseline, MMR with $\lambda = \frac{1}{3}$. This result is statistically significant and we show it using a permutation test. The *test statistic* is the difference in the percentage of times that the worker finds a matching movie in the lists generated by the best and second best performing methods. The *null hypothesis* is that all compared methods are equally good. Under this hypothesis, we permute the answers of the workers 10^6 times, generate an empirical distribution of the test statistic, and observe that the value of 13.5% or higher is less likely than 0.0001. So we reject the null hypothesis with $p < 0.0001$.

Secondly, the percentage of times that the worker considers the chosen movie to be a good genre-matching recommendation in the DUM list is 12.5% higher than in the list generated by the best performing baseline, MMR with $\lambda = \frac{1}{3}$. This result is statistically significant and we show it again using a permutation test. The *test*

²<http://www.imdb.com>

DUM	
The Shawshank Redemption	drama crime
The Dark Knight	drama thriller action crime
The Lord of the Rings 1	action adventure
Forrest Gump	drama romance
Back to the Future	comedy adventure
The Shining	drama horror
MMR ($\lambda = 1/3$)	
The Dark Knight	drama thriller action crime
Dr. Phibes Rises Again	comedy romance adventure horror
The Shawshank Redemption	drama crime
Pulp Fiction	thriller crime
Fight Club	drama
The Godfather	drama crime
MMR ($\lambda = 2/3$)	
The Dark Knight	drama thriller action crime
The Shawshank Redemption	drama crime
The Lord of the Rings 1	action adventure
Pulp Fiction	thriller crime
Fight Club	drama
The Godfather	drama crime
MMR ($\lambda = 0.99$)	
The Shawshank Redemption	drama crime
The Dark Knight	drama thriller action crime
Pulp Fiction	thriller crime
Fight Club	drama
The Godfather	drama crime
The Lord of the Rings 1	action adventure

Table 2: Four recommended lists in user study 1 where DUM outperforms MMR.

statistic is the difference in the percentage of times that the worker finds a good recommendation in the lists generated by the best and second best performing methods. The *null hypothesis* is that all compared methods are equally good. Under this hypothesis, we permute the answers of the workers 10^6 times, generate an empirical distribution of the test statistic, and observe that the value of 12.5% or higher is less likely than 0.001. So we reject the null hypothesis with $p < 0.001$. Overall, this user study shows that the diversity and the utility of recommendation lists generated by DUM are perceived superior to those of the lists generated by MMR.

We note that for all the methods compared in Table 1, the ratio between the percentage of times that the genre-matching movie is a good recommendation and that the matching movie is found is always between 0.92 and 0.94. This implies that if a matching movie found, it is very likely to be considered a good recommendation. We conjecture that this is due to the high popularity of movies in the ground set E , which practically guarantees the utility of the recommended movies and minimizes the differences between the compared methods.

In Table 2, we show a real-life example illustrating how DUM outperforms MMR. Here, DUM covers all the 8 movie genres by popular movies. These movies are well known and can be easily matched to any chosen target genre. However, MMR with $\lambda = 0.99$ assigns insufficient weight to diversity and therefore covers only five movie genres. The result is that this MMR list is unsuitable for users who like *Comedy*, *Romance*, or *Horror* movies. MMR with $\lambda = \frac{2}{3}$ has the same problem. On the other hand, MMR with $\lambda = \frac{1}{3}$ assigns too much weight to diversity and therefore covers four movie genres by a relatively unknown movie, “Phibes Rises Again”. These genres are not covered by any other movie in the list. The result is that this MMR list is likely to be of a low utility for users who like *Comedy*, *Romance*, *Adventure*, and *Horror* movies.

Instructions

Bob and Alice plan a vacation and can take several movies with them. **Bob loves drama movies** and **Alice loves romance movies**. Which movies should they take with them? Below are four lists of recommended movies. Please tell us what do you think about these lists. Choose the most appropriate judgment.

List 1

Forrest Gump
Titanic
Eternal Sunshine of the Spotless Mind
Slumdog Millionaire
The Shawshank Redemption
The Dark Knight
Fight Club
The Godfather

List 2

The Shawshank Redemption
The Dark Knight
Fight Club
The Godfather
The Lord of the Rings: The Return of the King
The Dark Knight Rises
Forrest Gump
Gladiator

List 3

The Shawshank Redemption
The Dark Knight
Fight Club
The Godfather
Forrest Gump
Titanic
Eternal Sunshine of the Spotless Mind
WALL-E

List 4

The Shawshank Redemption
The Dark Knight
Forrest Gump
Fight Club
The Godfather
Titanic
The Lord of the Rings: The Return of the King
Eternal Sunshine of the Spotless Mind

The list is good for

☐ Neither Bob nor Alice
☐ Bob who loves drama movies
☐ Alice who loves romance movies
☐ Both Bob and Alice

The list is good for

☐ Neither Bob nor Alice
☐ Bob who loves drama movies
☐ Alice who loves romance movies
☐ Both Bob and Alice

The list is good for

☐ Neither Bob nor Alice
☐ Bob who loves drama movies
☐ Alice who loves romance movies
☐ Both Bob and Alice

The list is good for

☐ Neither Bob nor Alice
☐ Bob who loves drama movies
☐ Alice who loves romance movies
☐ Both Bob and Alice

Figure 2: Our Mechanical Turk questionnaire in user study 2 for $t_1 = \text{Drama}$ and $t_2 = \text{Romance}$.

5.2 User Study 2

In the second study, we evaluate DUM on a specific problem of recommending a diverse set of movies that cover exactly two genres. We again compare DUM to three variants of MMR, which are parameterized by $\lambda \in \{\frac{1}{3}, \frac{2}{3}, 0.99\}$.

The compared methods are evaluated by MT workers. In each HIT, we ask the worker to consider a situation where Bob and Alice go for a vacation and can take several movies with them. Bob and Alice prefer two different movie genres. We generate four recommendation lists: one by DUM and three by MMR for different values of λ . For each list, we ask the worker to indicate whether the list is appropriate for both Bob and Alice, only for one of them, or for none of them. A screenshot of our MT questionnaire is shown in Figure 2.

Each HIT is associated with two movie genres, t_1 and t_2 , the preferences of Bob and Alice in the HIT. We generate three HITs for each pair of the 18 most frequent IMDb movie genres, so that the recommendation lists are evaluated $3 \frac{18 \times 17}{2} = 459$ times. Like in Section 5.1, the ground set E are 10k most frequently rated IMDb movies. The utility of movie e , $w(e)$, is the number of ratings assigned to this movie. The diversity function f is defined as in Lemma 2. The topics are $\mathcal{T} = \{t_1, t_2\}$ and $N_{t_1} = N_{t_2} = 4$. For this setting, DUM generates a list of at most 8 movies, at least 4 from each genre. The utility and diversity are normalized as in Section 5.1. In each HIT, the order of the recommendation lists is randomized and the length of the lists is determined as in Section 5.1.

Our HITs are completed by 57 *master* workers. Each worker is asked to complete at most 10 HITs. This guarantees that our HITs are completed by more than just a handful of workers. On average, a worker spends 57 seconds per HIT, i.e., 14 seconds to

Suitable for	DUM	MMR		
		$\lambda = \frac{1}{3}$	$\lambda = \frac{2}{3}$	$\lambda = 0.99$
Bob and Alice	74.51%	64.92%	58.39%	28.98%
Bob or Alice	23.53%	32.68%	39.43%	66.67%
Neither	1.96%	2.40%	2.18%	4.36%

Table 3: Comparison of DUM and MMR in user study 2. For each method, we report the percentage of times that the worker identifies the recommended list as suitable for both Bob and Alice; only for Bob or only for Alice; or for neither of them.

evaluate a list of up to 8 movies. In our analysis, we do not differentiate between suboptimal answers “Suitable only for Alice” and “Suitable only for Bob” and collapse the two into a single answer “Suitable for Alice or Bob”. The results of the second user study are presented in Table 3.

We observe that the workers consider the DUM list to be suitable for both Bob and Alice in 74.51% of cases. This is 9.6% higher than the best performing baseline, MMR with $\lambda = \frac{1}{3}$. This result is statistically significant and we show it using a permutation test. The *test statistic* is the difference in the percentage of times that the recommended lists, generated by the best and second best performing methods, are suitable for both Bob and Alice. The *null hypothesis* is that all compared methods are equally good. Under this hypothesis, we permute the answers of the workers 10^6 times, generate an empirical distribution of the test statistic, and observe that the value of 9.6% or higher is less likely than 0.0001. So we reject the null hypothesis with $p < 0.0001$.

Similarly to Section 5.1, our permutation test can be also interpreted as showing that our results are highly unlikely under the hypothesis that the workers are of low quality, the lists are rated randomly. This implies that our workers have reasonable expertise in evaluating our HITs.

In Table 4, we show another real-life example illustrating how DUM outperforms MMR for $t_1 = \text{Horror}$ and $t_2 = \text{Action}$. Here, DUM covers each movie genre by four most popular movies from that genre. However, MMR with $\lambda = 0.99$ assigns insufficient weight to diversity and therefore recommends only most popular items that happen to be *Action* movies. So the recommendation list is unsuitable for users who like *Horror* movies. MMR with $\lambda = \frac{2}{3}$ has a similar behavior and is strongly dominated by *Horror* movies. On the other hand, MMR with $\lambda = \frac{1}{3}$ assigns too much weight to diversity and therefore recommends many *Horror* movies that are also *Action* movies. These are less popular than the most popular *Horror* movies that are not *Action*. So the list is of a low utility for users who like *Horror* movies.

To sum up, DUM outperforms MMR in cases, where items from one topic have a higher utility than items from the other topic, and the items at the intersection of the two topics also have a low utility. While DUM recommends a mixture of high utility items from each topic, MMR either prefers items at the intersection of the topics, when the value of λ is low; or recommends high-utility items from the dominant topic only, when the value of λ is high.

5.3 Offline Evaluation

The main goal of the offline evaluation is to assess the performance of DUM under various conditions, such as recommendations across multiple users with their interest profiles defined based on different combinations of topics.

We use the *IM MovieLens* dataset [15] in the offline evaluation. The dataset consists of 1 million ratings on a 1-to-5 stars scale, assigned by about 6000 users to about 4000 movies. We remove

DUM	
The Dark Knight	action
The Lord of the Rings 1	action
The Matrix	action
Inception	action
The Shining	horror
Alien	horror
Psycho	horror
Shaun of the Dead	horror
MMR ($\lambda = 1/3$)	
Zombieland	horror action
From Dusk Till Dawn	horror action
Dawn of the Dead	horror action
Resident Evil	horror action
The Dark Knight	action
The Lord of the Rings 1	action
The Matrix	action
Inception	action
MMR ($\lambda = 2/3$)	
The Dark Knight	action
The Lord of the Rings 1	action
The Matrix	action
Inception	action
The Lord of the Rings 2	action
The Dark Knight Rises	action
The Lord of the Rings 3	action
The Shining	horror
MMR ($\lambda = 0.99$)	
The Dark Knight	action
The Lord of the Rings 1	action
The Matrix	action
Inception	action
The Lord of the Rings 2	action
The Dark Knight Rises	action
The Lord of the Rings 3	action
Avatar	action

Table 4: Four recommended lists in user study 2 where DUM outperforms MMR. The topics are $t_1 = \text{Horror}$ and $t_2 = \text{Action}$.

users having less than 300 ratings, so that for each user we have enough data to create a user profile and recommend movies. We end up with 1000 users and a total of 515k ratings.

Movies rated by each user are split randomly into the training and test set with the 2: 1 ratio. We use matrix factorization [22] to predict the rating of movies in the test set and feed the predicted ratings as the utility scores into the DUM and MMR methods. The split is performed three times for each user, and the reported results are based on the average of experiments conducted on these splits.

For each user, the training set is used for creating their interest profile, whereas the test set contains the recommendable movies (along with their actual and predicted utility). On average, a user profile is created based on 343 movies and the recommendation list is selected from a set of 171 candidates. These steps are carried out as follows:

User profile creation: There are 18 genres of movies in the dataset, and each movie belongs to one or more of these genres. For each user, we create a multinomial distribution over the popularity of genres of the movies rated by this user in the training data, assuming that users rated movies that they had watched. We sample 10 times from this distribution to create the user’s preference profile over genres, and normalize it so that the sum of the scores equals 1. For each user with the preference score r_t for genre t , we set $N_t = \lfloor r_t \times K \rfloor$ in (10), where K is the length of the recommendation list. That is, the coverage of each genre in the result list is proportional to the preference of the user for that genre.

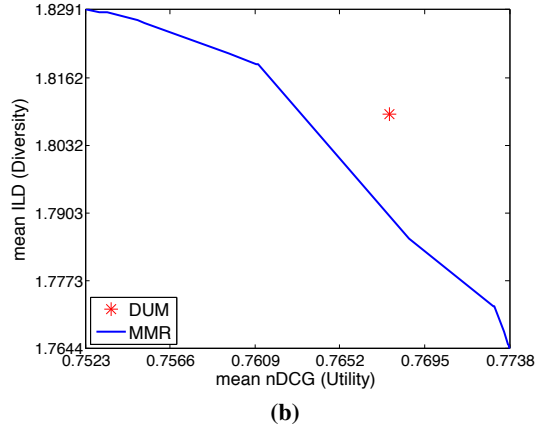
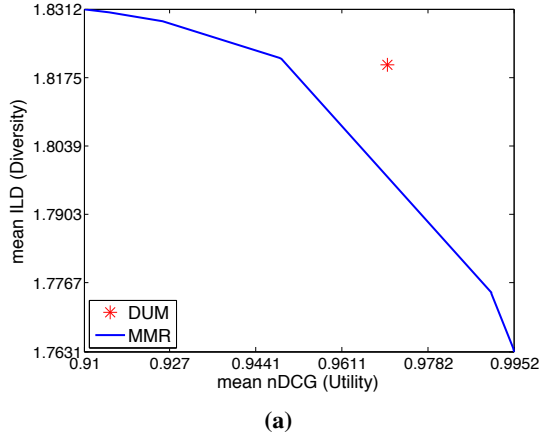


Figure 3: Performance of DUM in terms of diversity and utility is compared to the performance of MMR for all the settings of the parameter λ . (a) The actual rating of movies is the utility score, (b) The predicted rating of movies is the utility score.

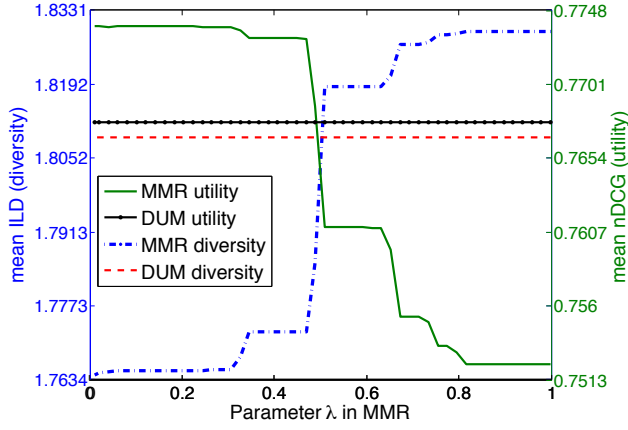


Figure 4: Tradeoff between the diversity and utility of the recommendation lists across all users for all the settings of the parameter λ in MMR. Diversity and utility scores achieved by DUM (independent of λ) are shown for comparison purposes.

Recommendation: Movies in the test data are used as the ground set E of recommendable movies, from which each diversification method finds the list of $K = 10$ movies to recommend to each user. The predicted utility of movies is used in the recommendation. The reason for using the predicted utility instead of the readily available movie ratings is to keep the evaluation as close as possible to real-world recommendation scenarios, where the utility of items is not known. When we evaluate the performance of the studied methods, we use the actual utility score, i.e., the rating assigned by a user to a test set movie.

5.3.1 Evaluation Metrics

We use three evaluation metrics to compare the performance of DUM to various settings of MMR: a diversity metric, a utility metric, and a compound metric that considers both diversity and utility. We chose these particular metrics due to two reasons. First, we wanted to evaluate the performance of our method with respect to diversity and utility individually (first two metrics), as well as in combination (third metric). Second, we wanted them to be different from the objective function of DUM in order to avoid any potential bias. Thus, the compound metric combines diversity and utility in

a different manner from what DUM does.

Intra-list distance (ILD) [24, 28] is a common metric that measures the diversity of a recommendation list as the average distance between pairs of recommended items. The dual of this measure is the intra-list similarity [31]. We use ILD to measure distance based diversity of a recommendation list in our experiment:

$$ILD = \frac{2}{|S|(|S| - 1)} \sum_{e \in S} \sum_{e' \in S} d(e, e') \quad (12)$$

where $d(e, e')$ measures the distance between two items e and e' in a list S . We choose the Euclidean distance between the genre vectors of two movies as the distance function d . Note that this metric is cardinaly different from the diversity function exploited by DUM, which is shown in (10).

Discounted cumulative gain (DCG) [12] measures the accumulated utility gain of items in the recommendations list from the top to the bottom, with the gain of each item s_k being discounted by its position k in the list:

$$DCG = \sum_{k=1}^{|S|} \frac{\mathbf{w}(s_k)}{\log(k+1)} \quad (13)$$

Here, $\mathbf{w}(s_k)$ is the utility of item s_k at rank k in the list. We estimate the utility of a movie for a user by the rating that the user assigned to the movie. We also use the normalized DCG (nDCG), which is in the range $[0, 1]$. nDCG is computed as $\text{nDCG} = \text{DCG}/\text{IDCG}$, where IDCG is the ideal gain achievable when all the listed items have the highest utility score.

Expected intra-list distance (EILD) [24] is a compound metric that combines utility and diversity. EILD measures the average intra-list distance (ILD) with respect to rank-sensitivity and utility:

$$EILD = \sum_{k=1}^{|S|} \sum_{k'=1}^{|S|} C_k \text{disc}(k) \text{rdisc}(k'|k) \mathbf{w}(s_k) \mathbf{w}(s_{k'}) d(s_k, s_{k'}) \quad (14)$$

where $\text{disc}(k) = 1/\log(k+1)$ is the discount function at rank k in the list and $\text{rdisc}(k'|k) = \text{disc}(\max(1, k' - k))$ is the relative rank discount. In order to avoid bias, we use the normalization constant proposed in [24] and set $C_k = \frac{1}{\bar{C}} / \sum_{k'=1}^{|S|} \text{disc}(k'|k) \mathbf{w}(s_{k'})$ where $\bar{C} = \sum_{k=1}^{|S|} \text{disc}(k)$.

We compute each of these metrics for every recommendation list

provided to a user. Then, we average them across the three runs for every user to compute user-based mean of the metric. This is performed for DUM and all settings of MMR, and the mean of each metric for each method is computed across all the users and reported.

5.3.2 Evaluation Results

Figure 3 shows the performance of DUM against MMR in terms of diversity and utility metrics. In Figure 3-a, the actual rating of movies in the test set is used as the utility of movies in the recommendation step, whereas in Figure 3-b the prediction produced by matrix factorization is used as the utility score. In both figures, MMR exhibits a trade-off between the values of mean ILD (as a measure of diversity) and mean nDCG (as a measure of utility). This trade-off is due to different values of the tuning parameter λ in different settings of MMR. For low values of λ the utility is prioritized, such that the diversity of lists generated by MMR is low, but the utility is high. An opposite situation is observed for high λ , when the diversity gets prioritized.

It can be seen that the performance of DUM with respect to both metrics is superior to any settings of MMR, regardless of the way the utility score is obtained. For instance, in Figures 3-b, DUM achieves nDCG of 0.767 (compared to the highest nDCG of 0.774 achieved by MMR for $\lambda = 0$) and ILD of 1.811 (compared to the highest ILD of 1.829 achieved by MMR for $\lambda = 1$). It should be highlighted that the utility and diversity cannot be optimized by MMR simultaneously, since they are achieved for different values of λ , while DUM competes with both of them at the same time. Also recall that DUM is parameter-free, and its superiority over MMR becomes clear.

Comparing Figures 3-a and 3-b, we observe that, as expected, the exact knowledge of the utility improves the performance of both DUM and MMR. However, this knowledge is unavailable to practical recommenders. Hence, we use the predicted utility values in the recommendation step in the following experiments, in order to mimic the conditions of a real-world recommendation scenario.

The trade-off between the diversity and utility objectives in MMR for various values of λ is visualized in Figure 4. As λ increases, the diversity of the list recommended by MMR increases whereas its utility decreases. It can be seen that $\lambda = 0.49$ is the operating point for MMR, where the utility and the diversity curves intersect. On the contrary, the utility and diversity of DUM are stable and both are above the operating point of MMR.

Another argument in favor of DUM is obtained through the EILD metric that combines diversity and utility. A comparison between DUM and all the possible settings of MMR with respect to EILD are plotted in Figure 5. It can be seen that DUM significantly outperforms MMR for low/moderate values of λ , which correspond to cases where utility is prioritized, or both utility and diversity are similarly important. MMR starts outperforming DUM for $\lambda > 0.65$, when the importance of diversity takes over, which may not be a favorite objective in real-world recommendations. This confirms the superiority of DUM in balancing the utility and diversity goals with no prior parameterization, compared to a method that explicitly targets the maximization of their weighted combination.

6. CONCLUSION

Much research in recommender systems has focused on the accuracy, but overlooked issues related to the composition of the recommendation lists. Increasing the diversity of the lists poses a trade-off to the utility, such that the problem of maximizing utility subject to diversity is an important challenge. In this work, we propose the diversity-weighted utility maximization (DUM) method and show that the problem can be cast as finding the maximum of a modular function on a polymatroid, which is known to have an op-

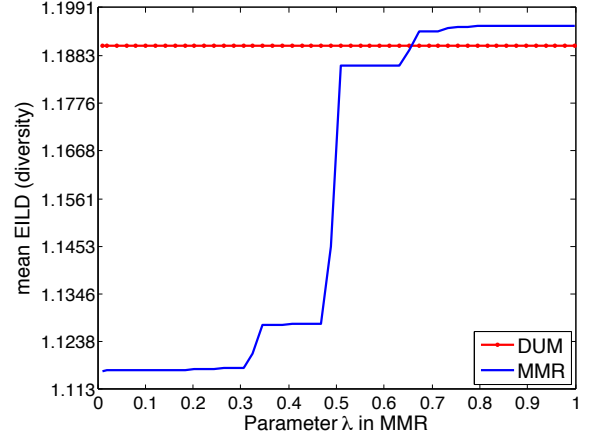


Figure 5: Performance of DUM in terms of expected diversity (w.r.t. utility and rank) is compared to the performance of MMR for all the settings of the parameter λ .

timal greedy solution. This parameter-free method guarantees that items in the recommendation list cover different aspects of user's taste, such that each aspect is covered by items with high utility.

We conduct two online user studies. The diversity and utility of DUM are evaluated in a movie recommendation scenario, and the perceived diversity of DUM is evaluated in a specific problem of recommending a diverse set of movies that cover exactly two genres. In both studies, we found that DUM outperforms baseline models that maximize a linear combination of utility and diversity. We also report an offline evaluation of DUM using a suite of diversity and utility metrics. The results show that DUM effectively balances the trade-off between diversity and utility: our method achieves performance comparable to the best performing baselines of diversity and utility, if executed individually. Moreover, a combined metric of diversity and utility shows the superiority of parameter-free DUM over the baseline methods that need to be parameterized.

Most diversification methods use MMR objective function, to linearly combine modular and submodular functions of utility and diversity, respectively. Our work is orthogonal to these methods in the sense that the DUM objective function maximizes a modular function subject to a submodular constraint. We demonstrate significant improvements over various settings of MMR, while we intend to conduct a more encompassing comparison with other variants MMR in the future. Another future direction is to account for the novelty of the recommended items [5, 6] with respect to prior consumption history of the user. This may be incorporated into the diversity function by considering, apart from the diversity contribution, also the novelty contribution of items in the list.

Another issue that deserves further investigation is the changes that need to be introduced in the diversity metric and in the tolerance for redundancy across different domains and applications. For instance, a metric of diversity applicable for news filtering may differ substantially from the metric we derived for the movie recommendation task in this work. Furthermore, user's tolerance for redundancy of news items that are in agreement with their own opinion may differ from their tolerance for redundancy of items having an opposite opinion. We intend to address these questions in our future works.

7. REFERENCES

- [1] R. Agrawal, S. Gollapudi, A. Halverson, and S. Ieong. Diversifying search results. In *Proceedings of the Second ACM International Conference on Web Search and Data Mining*, pages 5–14, 2009.
- [2] A. Ashkan, B. Kveton, S. Berkovsky, and Z. Wen. Diversified utility maximization for recommendations. In *Poster Proceedings of the 8th ACM Conference on Recommender Systems*, 2014.
- [3] G. Capannini, F. M. Nardini, R. Perego, and F. Silvestri. Efficient diversification of web search results. *Proceedings of the VLDB Endowment*, 4(7):451–459, 2011.
- [4] J. Carbonell and J. Goldstein. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 335–336, 1998.
- [5] P. Castells, S. Vargas, and J. Wang. Novelty and diversity metrics for recommender systems: choice, discovery and relevance. In *International Workshop on Diversity in Document Retrieval (DDR 2011) at the 33rd European Conference on Information Retrieval*, pages 29–36, 2011.
- [6] C. L. Clarke, M. Kolla, G. V. Cormack, O. Vechtomova, A. Ashkan, S. Büttcher, and I. MacKinnon. Novelty and diversity in information retrieval evaluation. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 659–666, 2008.
- [7] N. Craswell, O. Zoeter, M. Taylor, and B. Ramsey. An experimental comparison of click position-bias models. In *Proceedings of the 2008 International Conference on Web Search and Data Mining*, pages 87–94, 2008.
- [8] J. Edmonds. Submodular functions, matroids, and certain polyhedra. In *Combinatorial Structures and Their Applications: Proceedings of the Calgary International Conference on Combinatorial Structures and Their Applications*, pages 69–87. 1970.
- [9] S. Gollapudi and A. Sharma. An axiomatic approach for result diversification. In *Proceedings of the 18th WWW Conference*, pages 381–390, 2009.
- [10] M. Halvey, P. Punitha, D. Hannah, R. Villa, F. Hopfgartner, A. Goyal, and J. M. Jose. Diversity, assortment, dissimilarity, variety: A study of diversity measures using low level features for video retrieval. In *Advances in Information Retrieval*, pages 126–137. 2009.
- [11] D. Jannach, L. Lerche, F. Gedikli, and G. Bonnini. What recommenders recommend—an analysis of accuracy, popularity, and sales diversity effects. In *Proceedings of the 21st Conference on User Modeling, Adaptation, and Personalization*, pages 25–37. 2013.
- [12] K. Järvelin and J. Kekäläinen. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems (TOIS)*, 20(4):422–446, 2002.
- [13] Y. Koren and R. M. Bell. Advances in collaborative filtering. In *Recommender Systems Handbook*, pages 145–186. 2011.
- [14] A. Lad and Y. Yang. Learning to rank relevant and novel documents through user feedback. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 469–478, 2010.
- [15] S. Lam and J. Herlocker. MovieLens 1M Dataset. <http://www.grouplens.org/node/12>, 2014.
- [16] P. Lops, M. de Gemmis, and G. Semeraro. Content-based recommender systems: State of the art and trends. In *Recommender Systems Handbook*, pages 73–105. 2011.
- [17] S. M. McNee, J. Riedl, and J. A. Konstan. Being accurate is not enough: how accuracy metrics have hurt recommender systems. In *CHI’06 extended abstracts on Human factors in computing systems*, pages 1097–1101, 2006.
- [18] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing submodular set functions - I. *Mathematical Programming*, 14(1):265–294, 1978.
- [19] F. Radlinski, R. Kleinberg, and T. Joachims. Learning diverse rankings with multi-armed bandits. In *Proceedings of the 25th international conference on Machine learning*, pages 784–791, 2008.
- [20] F. Ricci, L. Rokach, B. Shapira, and P. B. Kantor, editors. *Recommender Systems Handbook*. Springer, 2011.
- [21] R. L. Santos, C. Macdonald, and I. Ounis. Exploiting query reformulations for web search result diversification. In *Proceedings of the 19th international conference on World wide web*, pages 881–890, 2010.
- [22] C. Thureau, K. Kersting, M. Wahabzada, and C. Bauckhage. Convex non-negative matrix factorization for massive datasets. *Knowledge and information systems*, 29(2):457–478, 2011.
- [23] D. Vallet and P. Castells. Personalized diversification of search results. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, pages 841–850, 2012.
- [24] S. Vargas and P. Castells. Rank and relevance in novelty and diversity metrics for recommender systems. In *Proceedings of the fifth ACM conference on Recommender systems*, pages 109–116, 2011.
- [25] S. Vargas, P. Castells, and D. Vallet. Explicit relevance models in intent-oriented information retrieval diversification. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, pages 75–84, 2012.
- [26] J. Yu, S. Mohan, D. P. Putthividhya, and W.-K. Wong. Latent dirichlet allocation based diversified retrieval for e-commerce search. In *Proceedings of the 7th ACM international conference on Web search and data mining*, pages 463–472, 2014.
- [27] Y. Yue and C. Guestrin. Linear submodular bandits and their application to diversified retrieval. In *Advances in Neural Information Processing Systems*, pages 2483–2491, 2011.
- [28] M. Zhang and N. Hurley. Avoiding monotony: improving the diversity of recommendation lists. In *Proceedings of the 2008 ACM conference on Recommender systems*, pages 123–130, 2008.
- [29] Y. C. Zhang, D. Ó. Séaghdha, D. Quercia, and T. Jambor. Auralist: introducing serendipity into music recommendation. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 13–22, 2012.
- [30] T. Zhou, Z. Kuscsik, J.-G. Liu, M. Medo, J. R. Wakeling, and Y.-C. Zhang. Solving the apparent diversity-accuracy dilemma of recommender systems. *Proceedings of the National Academy of Sciences*, 107(10):4511–4515, 2010.
- [31] C.-N. Ziegler, S. M. McNee, J. A. Konstan, and G. Lausen. Improving recommendation lists through topic diversification. In *Proceedings of the 14th international conference on World Wide Web*, pages 22–32, 2005.